

TITLE OF THE INVENTION

**TREE-BASED ORDERED MULTICASTING METHOD**

CROSS-REFERENCE TO RELATED APPLICATIONS

5           This application claims priority from U.S. provisional application serial number 60/244,405 filed on October 30, 2000, incorporated herein by reference.

STATEMENT REGARDING FEDERALLY SPONSORED RESEARCH  
OR DEVELOPMENT

10           This invention was made with Government support under Grant No. F19628-96-C-0338 awarded by the Air Force Office of Scientific Research (AFOSR).  
The Government has certain rights in this invention.

REFERENCE TO A COMPUTER PROGRAM APPENDIX

15           Not Applicable

NOTICE OF MATERIAL SUBJECT TO COPYRIGHT PROTECTION

20           A portion of the material in this patent document is subject to copyright protection under the copyright laws of the United States and of other countries. The owner of the copyright rights has no objection to the facsimile reproduction by anyone of the patent document or the patent disclosure, as it appears in the United States Patent and Trademark Office file or records, but otherwise reserves all copyright rights whatsoever.

The copyright owner does not hereby waive any of its rights to have this patent document maintained in secrecy, including without limitation its rights pursuant to 37 C.F.R. § 1.14.

## BACKGROUND OF THE INVENTION

### 5 1. Field of the Invention

This invention pertains generally to network multicast communication, and more particularly to ordering services for tree-based concurrent multicasting.

### 2. Description of the Background Art

100  
105  
110  
115  
120  
125  
130  
135  
140  
145  
150  
155  
160  
165  
170  
175  
180  
185  
190  
195  
200  
205  
210  
215  
220  
225  
230  
235  
240  
245  
250  
255  
260  
265  
270  
275  
280  
285  
290  
295  
300  
305  
310  
315  
320  
325  
330  
335  
340  
345  
350  
355  
360  
365  
370  
375  
380  
385  
390  
395  
400  
405  
410  
415  
420  
425  
430  
435  
440  
445  
450  
455  
460  
465  
470  
475  
480  
485  
490  
495  
500  
505  
510  
515  
520  
525  
530  
535  
540  
545  
550  
555  
560  
565  
570  
575  
580  
585  
590  
595  
600  
605  
610  
615  
620  
625  
630  
635  
640  
645  
650  
655  
660  
665  
670  
675  
680  
685  
690  
695  
700  
705  
710  
715  
720  
725  
730  
735  
740  
745  
750  
755  
760  
765  
770  
775  
780  
785  
790  
795  
800  
805  
810  
815  
820  
825  
830  
835  
840  
845  
850  
855  
860  
865  
870  
875  
880  
885  
890  
895  
900  
905  
910  
915  
920  
925  
930  
935  
940  
945  
950  
955  
960  
965  
970  
975  
980  
985  
990  
995  
1000  
1005  
1010  
1015  
1020  
1025  
1030  
1035  
1040  
1045  
1050  
1055  
1060  
1065  
1070  
1075  
1080  
1085  
1090  
1095  
1100  
1105  
1110  
1115  
1120  
1125  
1130  
1135  
1140  
1145  
1150  
1155  
1160  
1165  
1170  
1175  
1180  
1185  
1190  
1195  
1200  
1205  
1210  
1215  
1220  
1225  
1230  
1235  
1240  
1245  
1250  
1255  
1260  
1265  
1270  
1275  
1280  
1285  
1290  
1295  
1300  
1305  
1310  
1315  
1320  
1325  
1330  
1335  
1340  
1345  
1350  
1355  
1360  
1365  
1370  
1375  
1380  
1385  
1390  
1395  
1400  
1405  
1410  
1415  
1420  
1425  
1430  
1435  
1440  
1445  
1450  
1455  
1460  
1465  
1470  
1475  
1480  
1485  
1490  
1495  
1500  
1505  
1510  
1515  
1520  
1525  
1530  
1535  
1540  
1545  
1550  
1555  
1560  
1565  
1570  
1575  
1580  
1585  
1590  
1595  
1600  
1605  
1610  
1615  
1620  
1625  
1630  
1635  
1640  
1645  
1650  
1655  
1660  
1665  
1670  
1675  
1680  
1685  
1690  
1695  
1700  
1705  
1710  
1715  
1720  
1725  
1730  
1735  
1740  
1745  
1750  
1755  
1760  
1765  
1770  
1775  
1780  
1785  
1790  
1795  
1800  
1805  
1810  
1815  
1820  
1825  
1830  
1835  
1840  
1845  
1850  
1855  
1860  
1865  
1870  
1875  
1880  
1885  
1890  
1895  
1900  
1905  
1910  
1915  
1920  
1925  
1930  
1935  
1940  
1945  
1950  
1955  
1960  
1965  
1970  
1975  
1980  
1985  
1990  
1995  
2000  
2005  
2010  
2015  
2020  
2025  
2030  
2035  
2040  
2045  
2050  
2055  
2060  
2065  
2070  
2075  
2080  
2085  
2090  
2095  
2100  
2105  
2110  
2115  
2120  
2125  
2130  
2135  
2140  
2145  
2150  
2155  
2160  
2165  
2170  
2175  
2180  
2185  
2190  
2195  
2200  
2205  
2210  
2215  
2220  
2225  
2230  
2235  
2240  
2245  
2250  
2255  
2260  
2265  
2270  
2275  
2280  
2285  
2290  
2295  
2300  
2305  
2310  
2315  
2320  
2325  
2330  
2335  
2340  
2345  
2350  
2355  
2360  
2365  
2370  
2375  
2380  
2385  
2390  
2395  
2400  
2405  
2410  
2415  
2420  
2425  
2430  
2435  
2440  
2445  
2450  
2455  
2460  
2465  
2470  
2475  
2480  
2485  
2490  
2495  
2500  
2505  
2510  
2515  
2520  
2525  
2530  
2535  
2540  
2545  
2550  
2555  
2560  
2565  
2570  
2575  
2580  
2585  
2590  
2595  
2600  
2605  
2610  
2615  
2620  
2625  
2630  
2635  
2640  
2645  
2650  
2655  
2660  
2665  
2670  
2675  
2680  
2685  
2690  
2695  
2700  
2705  
2710  
2715  
2720  
2725  
2730  
2735  
2740  
2745  
2750  
2755  
2760  
2765  
2770  
2775  
2780  
2785  
2790  
2795  
2800  
2805  
2810  
2815  
2820  
2825  
2830  
2835  
2840  
2845  
2850  
2855  
2860  
2865  
2870  
2875  
2880  
2885  
2890  
2895  
2900  
2905  
2910  
2915  
2920  
2925  
2930  
2935  
2940  
2945  
2950  
2955  
2960  
2965  
2970  
2975  
2980  
2985  
2990  
2995  
3000  
3005  
3010  
3015  
3020  
3025  
3030  
3035  
3040  
3045  
3050  
3055  
3060  
3065  
3070  
3075  
3080  
3085  
3090  
3095  
3100  
3105  
3110  
3115  
3120  
3125  
3130  
3135  
3140  
3145  
3150  
3155  
3160  
3165  
3170  
3175  
3180  
3185  
3190  
3195  
3200  
3205  
3210  
3215  
3220  
3225  
3230  
3235  
3240  
3245  
3250  
3255  
3260  
3265  
3270  
3275  
3280  
3285  
3290  
3295  
3300  
3305  
3310  
3315  
3320  
3325  
3330  
3335  
3340  
3345  
3350  
3355  
3360  
3365  
3370  
3375  
3380  
3385  
3390  
3395  
3400  
3405  
3410  
3415  
3420  
3425  
3430  
3435  
3440  
3445  
3450  
3455  
3460  
3465  
3470  
3475  
3480  
3485  
3490  
3495  
3500  
3505  
3510  
3515  
3520  
3525  
3530  
3535  
3540  
3545  
3550  
3555  
3560  
3565  
3570  
3575  
3580  
3585  
3590  
3595  
3600  
3605  
3610  
3615  
3620  
3625  
3630  
3635  
3640  
3645  
3650  
3655  
3660  
3665  
3670  
3675  
3680  
3685  
3690  
3695  
3700  
3705  
3710  
3715  
3720  
3725  
3730  
3735  
3740  
3745  
3750  
3755  
3760  
3765  
3770  
3775  
3780  
3785  
3790  
3795  
3800  
3805  
3810  
3815  
3820  
3825  
3830  
3835  
3840  
3845  
3850  
3855  
3860  
3865  
3870  
3875  
3880  
3885  
3890  
3895  
3900  
3905  
3910  
3915  
3920  
3925  
3930  
3935  
3940  
3945  
3950  
3955  
3960  
3965  
3970  
3975  
3980  
3985  
3990  
3995  
4000  
4005  
4010  
4015  
4020  
4025  
4030  
4035  
4040  
4045  
4050  
4055  
4060  
4065  
4070  
4075  
4080  
4085  
4090  
4095  
4100  
4105  
4110  
4115  
4120  
4125  
4130  
4135  
4140  
4145  
4150  
4155  
4160  
4165  
4170  
4175  
4180  
4185  
4190  
4195  
4200  
4205  
4210  
4215  
4220  
4225  
4230  
4235  
4240  
4245  
4250  
4255  
4260  
4265  
4270  
4275  
4280  
4285  
4290  
4295  
4300  
4305  
4310  
4315  
4320  
4325  
4330  
4335  
4340  
4345  
4350  
4355  
4360  
4365  
4370  
4375  
4380  
4385  
4390  
4395  
4400  
4405  
4410  
4415  
4420  
4425  
4430  
4435  
4440  
4445  
4450  
4455  
4460  
4465  
4470  
4475  
4480  
4485  
4490  
4495  
4500  
4505  
4510  
4515  
4520  
4525  
4530  
4535  
4540  
4545  
4550  
4555  
4560  
4565  
4570  
4575  
4580  
4585  
4590  
4595  
4600  
4605  
4610  
4615  
4620  
4625  
4630  
4635  
4640  
4645  
4650  
4655  
4660  
4665  
4670  
4675  
4680  
4685  
4690  
4695  
4700  
4705  
4710  
4715  
4720  
4725  
4730  
4735  
4740  
4745  
4750  
4755  
4760  
4765  
4770  
4775  
4780  
4785  
4790  
4795  
4800  
4805  
4810  
4815  
4820  
4825  
4830  
4835  
4840  
4845  
4850  
4855  
4860  
4865  
4870  
4875  
4880  
4885  
4890  
4895  
4900  
4905  
4910  
4915  
4920  
4925  
4930  
4935  
4940  
4945  
4950  
4955  
4960  
4965  
4970  
4975  
4980  
4985  
4990  
4995  
5000  
5005  
5010  
5015  
5020  
5025  
5030  
5035  
5040  
5045  
5050  
5055  
5060  
5065  
5070  
5075  
5080  
5085  
5090  
5095  
5100  
5105  
5110  
5115  
5120  
5125  
5130  
5135  
5140  
5145  
5150  
5155  
5160  
5165  
5170  
5175  
5180  
5185  
5190  
5195  
5200  
5205  
5210  
5215  
5220  
5225  
5230  
5235  
5240  
5245  
5250  
5255  
5260  
5265  
5270  
5275  
5280  
5285  
5290  
5295  
5300  
5305  
5310  
5315  
5320  
5325  
5330  
5335  
5340  
5345  
5350  
5355  
5360  
5365  
5370  
5375  
5380  
5385  
5390  
5395  
5400  
5405  
5410  
5415  
5420  
5425  
5430  
5435  
5440  
5445  
5450  
5455  
5460  
5465  
5470  
5475  
5480  
5485  
5490  
5495  
5500  
5505  
5510  
5515  
5520  
5525  
5530  
5535  
5540  
5545  
5550  
5555  
5560  
5565  
5570  
5575  
5580  
5585  
5590  
5595  
5600  
5605  
5610  
5615  
5620  
5625  
5630  
5635  
5640  
5645  
5650  
5655  
5660  
5665  
5670  
5675  
5680  
5685  
5690  
5695  
5700  
5705  
5710  
5715  
5720  
5725  
5730  
5735  
5740  
5745  
5750  
5755  
5760  
5765  
5770  
5775  
5780  
5785  
5790  
5795  
5800  
5805  
5810  
5815  
5820  
5825  
5830  
5835  
5840  
5845  
5850  
5855  
5860  
5865  
5870  
5875  
5880  
5885  
5890  
5895  
5900  
5905  
5910  
5915  
5920  
5925  
5930  
5935  
5940  
5945  
5950  
5955  
5960  
5965  
5970  
5975  
5980  
5985  
5990  
5995  
6000  
6005  
6010  
6015  
6020  
6025  
6030  
6035  
6040  
6045  
6050  
6055  
6060  
6065  
6070  
6075  
6080  
6085  
6090  
6095  
6100  
6105  
6110  
6115  
6120  
6125  
6130  
6135  
6140  
6145  
6150  
6155  
6160  
6165  
6170  
6175  
6180  
6185  
6190  
6195  
6200  
6205  
6210  
6215  
6220  
6225  
6230  
6235  
6240  
6245  
6250  
6255  
6260  
6265  
6270  
6275  
6280  
6285  
6290  
6295  
6300  
6305  
6310  
6315  
6320  
6325  
6330  
6335  
6340  
6345  
6350  
6355  
6360  
6365  
6370  
6375  
6380  
6385  
6390  
6395  
6400  
6405  
6410  
6415  
6420  
6425  
6430  
6435  
6440  
6445  
6450  
6455  
6460  
6465  
6470  
6475  
6480  
6485  
6490  
6495  
6500  
6505  
6510  
6515  
6520  
6525  
6530  
6535  
6540  
6545  
6550  
6555  
6560  
6565  
6570  
6575  
6580  
6585  
6590  
6595  
6600  
6605  
6610  
6615  
6620  
6625  
6630  
6635  
6640  
6645  
6650  
6655  
6660  
6665  
6670  
6675  
6680  
6685  
6690  
6695  
6700  
6705  
6710  
6715  
6720  
6725  
6730  
6735  
6740  
6745  
6750  
6755  
6760  
6765  
6770  
6775  
6780  
6785  
6790  
6795  
6800  
6805  
6810  
6815  
6820  
6825  
6830  
6835  
6840  
6845  
6850  
6855  
6860  
6865  
6870  
6875  
6880  
6885  
6890  
6895  
6900  
6905  
6910  
6915  
6920  
6925  
6930  
6935  
6940  
6945  
6950  
6955  
6960  
6965  
6970  
6975  
6980  
6985  
6990  
6995  
7000  
7005  
7010  
7015  
7020  
7025  
7030  
7035  
7040  
7045  
7050  
7055  
7060  
7065  
7070  
7075  
7080  
7085  
7090  
7095  
7100  
7105  
7110  
7115  
7120  
7125  
7130  
7135  
7140  
7145  
7150  
7155  
7160  
7165  
7170  
7175  
7180  
7185  
7190  
7195  
7200  
7205  
7210  
7215  
7220  
7225  
7230  
7235  
7240  
7245  
7250  
7255  
7260  
7265  
7270  
7275  
7280  
7285  
7290  
7295  
7300  
7305  
7310  
7315  
7320  
7325  
7330  
7335  
7340  
7345  
7350  
7355  
7360  
7365  
7370  
7375  
7380  
7385  
7390  
7395  
7400  
7405  
7410  
7415  
7420  
7425  
7430  
7435  
7440  
7445  
7450  
7455  
7460  
7465  
7470  
7475  
7480  
7485  
7490  
7495  
7500  
7505  
7510  
7515  
7520  
7525  
7530  
7535  
7540  
7545  
7550  
7555  
7560  
7565  
7570  
7575  
7580  
7585  
7590  
7595  
7600  
7605  
7610  
7615  
7620  
7625  
7630  
7635  
7640  
7645  
7650  
7655  
7660  
7665  
7670  
7675  
7680  
7685  
7690  
7695  
7700  
7705  
7710  
7715  
7720  
7725  
7730  
7735  
7740  
7745  
7750  
7755  
7760  
7765  
7770  
7775  
7780  
7785  
7790  
7795  
7800  
7805  
7810  
7815  
7820  
7825  
7830  
7835  
7840  
7845  
7850  
7855  
7860  
7865  
7870  
7875  
7880  
7885  
7890  
7895  
7900  
7905  
7910  
7915  
7920  
7925  
7930  
7935  
7940  
7945  
7950  
7955  
7960  
7965  
7970  
7975  
7980  
7985  
7990  
7995  
8000  
8005  
8010  
8015  
8020  
8025  
8030  
8035  
8040  
8045  
8050  
8055  
8060  
8065  
8070  
8075  
8080  
8085  
8090  
8095  
8100  
8105  
8110  
8115  
8120  
8125  
8130  
8135  
8140  
8145  
8150  
8155  
8160  
8165  
8170  
8175  
8180  
8185  
8190  
8195  
8200  
8205  
8210  
8215  
8220  
8225  
8230  
8235  
8240  
8245  
8250  
8255  
8260  
8265  
8270  
8275  
8280  
8285  
8290  
8295  
8300  
8305  
8310  
8315  
8320  
8325  
8330  
8335  
8340  
8345  
8350  
8355  
8360  
8365  
8370  
8375  
8380  
8385  
8390  
8395  
8400  
8405  
8410  
8415  
8420  
8425  
8430  
8435  
8440  
8445  
8450  
8455  
8460  
8465  
8470  
8475  
8480  
8485  
8490  
8495  
8500  
8505  
8510  
8515  
8520  
8525  
8530  
8535  
8540  
8545  
8550  
8555  
8560  
8565  
8570  
8575  
8580  
8585  
8590  
8595  
8600  
8605  
8610  
8615  
8620  
8625  
8630  
8635  
8640  
8645  
8650  
8655  
8660  
8665  
8670  
8675  
8680  
8685  
8690  
8695  
8700  
8705  
8710  
8715  
8720  
8725  
8730  
8735  
8740  
8745  
8750  
8755  
8760  
8765  
8770  
8775  
8780  
8785  
8790  
8795  
8800  
8805  
8810  
8815  
8820  
8825  
8830  
8835  
8840  
8845  
8850  
8855  
8860  
8865  
8870  
8875  
8880  
8885  
8890  
8895  
8900  
8905  
8910  
8915  
8920  
8925  
8930  
8935  
8940  
8945  
8950  
8955  
8960  
8965  
8970  
8975  
8980  
8985  
8990  
8995  
9000  
9005  
9010  
9015  
9020  
9025  
9030  
9035  
9040  
9045  
9050  
9055  
9060  
9065  
9070  
9075  
9080  
9085  
9090  
9095  
9100  
9105  
9110  
9115  
9120  
9125  
9130  
9135  
9140  
9145  
9150  
9155  
9160  
9165  
9170  
9175  
9180  
9185  
9190  
9195  
9200  
9205  
9210  
9215  
9220  
9225  
9230  
9235  
9240  
9245  
9250  
9255  
9260  
9265  
9270  
9275  
9280  
9285  
9290  
9295  
9300  
9305  
9310  
9315  
9320  
9325  
9330  
9335  
9340  
9345  
9350  
9355  
9360  
9365  
9370  
9375  
9380  
9385  
9390  
9395  
9400  
9405  
9410  
9415  
9420  
9425  
9430  
9435  
9440  
9445  
9450  
9455  
9460  
9465  
9470  
9475  
9480  
9485  
9490  
9495  
9500  
9505  
9510  
9515  
9520  
9525  
9530  
9535  
9540  
9545  
9550  
9555  
9560  
9565  
9570  
9575  
9580  
9585  
9590  
9595  
9600  
9605  
9610  
9615  
9620  
9625  
9630  
9635  
9640  
9645  
9650  
9655  
9660  
9665  
9670  
9675  
9680  
9685  
9690  
9695  
9700  
9705  
9710  
9715  
9720  
9725  
9730  
9735  
9740  
9745  
9750  
9755  
9760  
9765  
9770  
9775  
9780  
9785  
9790  
9795  
9800  
9805  
9810  
9815  
9820  
9825  
9830  
9835  
9840  
9845  
9850  
9855  
9860  
9865  
9870  
9875  
9880  
9885  
9890  
9895  
9900  
9905  
9910  
9915  
9920  
9925  
9930  
9935  
9940  
9945  
9950  
9955  
9960  
9965  
9970  
9975  
9980  
9985  
9990  
9995  
10000  
10005  
10010  
10015  
10020  
10025  
10030  
10035  
10040  
10045  
10050  
10055  
10060  
10065  
10070  
10075  
10080  
10085  
10090  
10095  
10100  
10105  
10110  
10115  
10120  
10125  
10130  
10135  
10140  
10145  
10150  
10155  
10160  
10165  
10170  
10175  
10180  
10185  
10190  
10195  
10200  
10205  
10210  
10215  
10220  
10225  
10230  
10235  
10240  
10245  
10250  
10255  
10260  
10265  
10270  
10275  
10

or deferred to the application layer. A large body of work in the field of total and causal ordering for multicast messages is centered around fault tolerance or consistency issues in distributed systems.

Therefore, a need exists for a method of ordered multicasting that operates directly on reliable multicast trees to provide increased scalability, efficiency, and practicality. The present invention satisfies those needs, as well as others, and overcomes the deficiencies of previously developed multicast protocols.

### BRIEF SUMMARY OF THE INVENTION

The present invention comprises a solution for message ordering services integrated with a tree-based, concurrent, reliable multicast. Multicasting is essential for efficient one-to-many communications in a computer network. The Internet infrastructure and applications are increasingly being adapted to multicasting and require reliability and effective ordering of message transmissions. While reliability has been extensively researched in recent years, a solution for integrated ordered delivery over the most common delivery geometries (trees) within the Internet has been lacking, and is provided within this present invention.

According to an aspect of the invention, ordering is performed on a tree, instead of a ring, as proposed by prior work on reliable multicast protocols. The ordering process is performed on a mirror copy of an underlying shared multicast tree and supports ordering of messages from rapidly changing sources, for overlapping receiver groups and for anonymous hosts.

Ordering can be deployed more practically as a middleware component for any application needing ordered delivery, as opposed to requiring each application to provide its own, independent, ordering service. Ordering within the present invention is distributed among many nodes across the tree and thereby achieves improved scalability and efficiency.

By way of further example, and not of limitation, the invention provides ordering of messages for applications using IP multicasting within the Internet. A novel taxonomy of ordered broadcast and multicast solutions and a basic comparison of message complexities indicates that using the underlying infrastructure of trees predominant in current IP-multicasting solutions achieves the same or better efficiency in comparison with previous approaches. Support for ordering below the application level allows more rapid design and deployment of applications depending on ordered multicasting. Previous work on reliable multicasting indicated that shared trees provided the most efficient infrastructure for reliable data dissemination. Shared trees allow for concurrent usage of the same tree geometry by multiple sources disseminating data to different groups on the tree. The tree-based ordered multicasting (TOM) protocol of the present invention adds total ordering of packets to concurrent reliable multicast, wherein the ordering operation is distributed across the nodes within the network. A number of features are provided within the TOM to facilitate the ordering operation. A mirror copy of a logical tree geometry is utilized to provide concurrent, reliable multicasting as an infrastructure for ordering. Aggregation of ordering primitives is performed to minimize control traffic among nodes, in resemblance to a two-phase

ordering protocol, however, it is deployed across the tree. Aggregation entails the ordering and combination of messages destined for the same receivers, performed at hosts on the delivery path. TOM utilizes address extensions assigned to hosts for self-routing of messages and dynamic distribution of the ordering processing load. By using the address extensions, TOM also supports total ordering of messages for anonymous and overlapping receiver groups in shared trees, and can be extended to support causal and atomic ordered multicast. The use of causal and atomic multicast can also be supported with minor changes in the protocol delivery semantics. The ordered multicast, as described and specified with the TOM protocol, can be implemented in either software or hardware.

An object of the invention is to provide ordered multicasting for tree-based multicasting networks.

Another object of the invention is to provide ordered multicasting which employs distributed ordering responsibilities across the tree.

Another object of the invention is to provide for ordered multicasting with improved scalability, resiliency, and efficiency, of the concurrent transmissions.

Another object of the invention is to provide ordered multicasting with integrated reliability provisions and ordering in the same topology and delivery process.

Another object of the invention is to provide ordered multicasting in which extra computations and maintenance of a propagation graph are not necessary.

Another object of the invention is to provide ordered multicasting that allows ordered concurrent transmissions from rapidly changing sources on the same tree.

Another object of the invention is to provide ordered multicasting in which address extensions allow dynamic election of any node on the tree to order messages destined for the same group.

Another object of the invention is to provide ordered multicasting in which the address extensions support ordered delivery to anonymous hosts and overlapping receiver groups in shared trees.

Further objects and advantages of the invention will be brought out in the following portions of the specification, wherein the detailed description is for the purpose of fully disclosing preferred embodiments of the invention without placing limitations thereon.

#### BRIEF DESCRIPTION OF THE DRAWINGS

The invention will be more fully understood by reference to the following drawings which are for illustrative purposes only:

FIG. 1 is a protocol stack diagram of ordered multicasting according to an embodiment of the present invention, as shown as middleware within the host software.

FIG. 2 is a flowchart of multicasting operation according to an embodiment of the TOM protocol of the present invention.

FIG. 3 is a topology diagram upon which the operation of the TOM protocol according to an embodiment of the present invention is exemplified.

FIG. 4 is a pseudocode listing of TOM procedures according to aspects of the present invention, showing send, receive, and casting procedures.

FIG. 5 is a tree-diagram showing the classifications of ordering paradigms wherein the TOM protocol according to the present invention, showing the TOM protocol classified as a tree-based geometry.

FIG. 6 is a graph of multicast message costs which compares a number of protocols, including the TOM protocol according to an embodiment of the present invention.

## DETAILED DESCRIPTION OF THE INVENTION

For illustrative purposes the present invention will be described with reference to FIG. 1 through FIG. 6. It will be appreciated that the apparatus may vary as to configuration and as to details of the parts, and that the method may vary as to the specific steps and sequence, without departing from the basic concepts as disclosed herein.

### 1. Introduction

IP multicast communication generalizes the point-to-point and broadcast communication model to multipoint dissemination of messages. A source is required to transmit a single stream of packets to the network interface whereupon those packets are transparently replicated along their transmission paths to the receivers. This form of communication is indispensable for networked applications with high-volume data transfer, such as distributed software updates, news casts, video-on-demand, and interactive applications which include distributed simulations and telecollaboration systems. Data handled by these applications fall into two categories, continuous media streams and non-real-time data. Real-time data delivery, such as utilized for delivering

video or audio streams, is typically best-effort and unordered, but must observe deadlines to be useful for an application. Non-real time packets carry discrete data, and may require reliable, ordered, delivery based on the application semantics.

Changes in datagram routing or transmission errors may cause packets to arrive at their destination out of sequence. Disordered delivery of packets in a distributed application may result in different views of the group state at end hosts. Ordering of messages compensates for the lack of a global system state and the effects of asynchrony, unpredictable network delay, and disparities in host processing in distributed communication, while its use warrants that destination processes observe the same order of reception of messages. The ordering of messages is complemented by reliability and atomicity. Reliability guarantees that messages eventually arrive correctly at their destinations, while atomicity guarantees that a message is received by all members of a multicast group or none.

Consider a distributed interactive simulation with many moving, interacting entities, wherein a message  $m_1$  is reliably multicast from source  $s_1$  to receiver group  $Rec_1$ , and  $m_2$  is reliably multicast from  $s_2$  to  $Rec_2$ . A host which belongs to  $Rec_1$  may receive message  $m_1$  before  $m_2$ , while another host belonging to both groups may receive the messages in the opposite order. Correct operation of the simulation system requires not only that the input stream is equivalent for all replicas, but all input events have to be delivered to the replicated instances of shared applications in the same order. An ordering protocol must intercept, or preferably be integrated within, the



delivery process to guarantee the described consistency.

The majority of current reliable multicasting solutions lack associated ordering services. In a performance comparison of such protocols, entailing both sender and receiver initiated protocols, ring or tree-based protocols, and tree protocols with negative acknowledgments and periodic polling, it was determined that the latter protocol type was the most scalable and efficient approach known to date among deployable systems. Based on these observations, our objective is to examine how ordering services can be integrated with reliable multicasting, in particular with tree-based protocols, preserving scalability and efficiency. The present invention provides a solution for this problem using staggered ordering of messages on their delivery paths from sources to receivers in the reliable multicast tree, which is also used for logical connectivity between hosts for the purpose of error recovery. In contrast to earlier work, the protocol of the present invention does not require construction of a separate logical propagation graph or global clock synchronization, and ordering is distributed across nodes on the delivery paths between sources and receivers in the multicast tree.

## 2. System Model and Assumptions

The present network model  $= (H, C)$  consists of a set of  $k$  hosts  $H$  and communication links  $C$ , communicating via message passing in the absence of physical clock synchronization. A host is equated with the processes running on it. A multicast group is a set of  $k$  hosts in a network of  $H$  hosts, which is addressable collectively by a unique group address.

Message dissemination is assumed to be genuine multicast, such as wherein a source sends a message  $m$  once to the network interface in a multicast enabled backbone, which replicates  $m$  at multicast enabled routers on its path to  $r \leq n$  receivers. This stands in contrast to most prior work on ordered multicasting which assumes either unicast, where a message must be sent  $r$  times from a source to the network interface to reach  $r < n$  receivers, or broadcast, wherein all  $n$  hosts in the network are addressed and designated receivers must filter out messages targeted at them.

Four cases of group connectivity can be observed, (1) from a single source  $s$  to a single group  $g$ , denoted as  $(s, g)$ ; or (2) to multiple groups  $G$ ,  $(s, G)$ , or from multiple sources  $S$  to; (3) a single group,  $(S, g)$ , or (4) to multiple groups,  $(S, G)$ . Cases (1) and (2) have a trivial solution wherein sequence numbers fixing the ordering relation are added to outgoing messages at the source and are delivered in that order at the destinations. Cases (3) and (4), however, are more difficult to implement, because sending messages from one host is independent from the other hosts, whereas reception of the same messages may be interdependent and destination groups may overlap.

The present methods are directed toward totally ordered multicasts from multiple sources to multiple receivers or receiver groups. It is assumed that hosts do not fail and that network partitions do not occur. Overlapping groups are also considered in relation to the present protocol, as these were a focal point in previous work on ordered

multicast. Hosts contained within the intersection of two overlapping multicast groups should receive a message only once if the message is sent to both groups.

In total order, two messages  $m_1$  and  $m_2$  are sent to a receiver set  $Rec$  in the same relative order. For example, if two sources,  $A$  and  $B$ , send messages  $m_1$  and  $m_2$  to receiver groups  $G_1$  and  $G_2$ , respectively, then hosts in both groups, in particular in the intersection  $G_1 \cap G_2$ , should receive both messages either in the order  $(m_1, m_2)$ , or  $(m_2, m_1)$ . Atomic order demands that either all or none of the hosts in  $Rec$  receive the messages. A weaker notion of total order is causal order, based on Lamport's "happened before" relation. While a causal precedence relation between two messages preserves their sending order at delivery time, messages without causal linkage may still be delivered to different hosts in different order. Logical point-to-point channels between any pair of hosts are assumed to be FIFO to prevent an earlier message by the same process from being overtaken during delivery by a later message. If not provided by the network layer, FIFO-delivery over non-FIFO channels can be implemented by having the source process add a sequence number to its messages and let destinations deliver according to such sequence numbers.

Finally, it is assumed that a reliable, unordered multicast protocol is running at every host providing reliable delivery of a message to all operational hosts in a target multicast group. Ordered multicast should be *host minimal*, wherein no other hosts are affected by multicasting of the message other than the source and receivers, and *message minimal*, wherein the message size is a function of the size of the receiver set



know the identity of all receivers in the multicast group. However, the paths from sources to receivers may be suboptimal.

Although a reliable multicast protocol should be utilized with the present ordering mechanism, it is unimportant for the present description to specify a particular multicast protocol. The use of source-based or shared dissemination is also not crucial, however, the present invention will exemplify the operation of TOM to provide total order in a shared tree. An important concept in TOM is to multicast a message from a source to a receiver set combined with sending ordering information for the message, such as sequence numbers or time-stamps, to a common node on the tree which has been elected as the ordering node for this receiver set, or multicast group. The ordering node is responsible for sequencing the messages assigned to it and multicasts binding sequence numbers for final delivery to the receiver set, wherein the pending messages are to be delivered. TOM can be deployed in the form of an API accessible to applications with ordering needs.

### 3.1. Data Structures

A host in the multicast tree is either a source node (SN), an extra node (EN), a primary node (PN), an ordering node (ON), or a receiver node (RN). Since every host in the multicast session runs the ordering protocol, roles are assumed on-the-fly and no dedicated hardware is needed. The source node, SN, emit messages to one or more multicast groups in a session. Each extra node, EN, is a node that is not a member of the receiver set for a message, relaying messages upward or downward in the tree without participation in the ordering process. Primary nodes, PNs, are hosts on the

upward ordering path from source node, SN, to ordering node, ON, aggregating control messages in local order and forwarding revised sequence numbers up in the tree. The ordering node, ON, is the sequencer node for a message, gathering sequence number bids set *en route* by primary node, PN, deciding on a globally valid number, and multicasting the message to the receiver set with a final and binding sequence number directive. Sources can be ordering nodes, ONs, as well. Receiver nodes, RNs, are recipients of message which are delivering them according to an ordering-node, ON, sanctioned sequence number. Nodes can be source nodes, SN, for their own messages and assume all other roles for other messages. Edges within the acknowledgment-tree point from child nodes to their parent nodes.

A TOM message  $m = (m^h, m^b)$  consists of a control header  $m^h$  and body  $m^b$ , with  $m^h = (SN\_id, Rec, Seq\#, ts, of)$  where  $SN\_id$  is the source node identifier,  $Rec$  is the target receiver set, which is either a multicast group, or a collection of individual node identifiers;  $Seq\#$  is the sequence number used for ordering,  $ts$  is an optional time-stamp for ordering in response to timing information at the nodes, and  $of$  is the ordering flag indicating that a binding sequence number for the message has been set, while  $m^b$  contains the actual data stream.

Each node maintains two message windows for ordering, with a window for unordered messages ( $uw$ ), which have been received but whose delivery is pending; and an ordered messages window ( $ow$ ) for messages, which are correctly ordered and can be delivered to local processes. The sizes of these buffers are limited by the

number of hosts in the largest multicast group known at the time of buffer allocation.

Each host programs its local network interface to subscribe to multicast packets on the same local network, or to receive packets from routers based on IGMP information .

### 3.2 Operation of TOM

FIG. 2 illustrates the general operation of the TOM protocol for ordering multicast messages according to four steps: first, a message multicast from each source node, SN, to receivers as shown by block 100; next a control message unicast from a source node, SN, across a primary node, PN, to the ordering node, ON, for the designated multicast group or transmission as per block 102, where primary node, PN, aggregates messages from their subtrees and hence staggers the ordering process upward within the tree; then, determination of a binding sequence number for this message and a multicast to the receiver group as shown in block 144; and finally, the delivery of messages at end hosts according to the agreed-upon sequence numbers as per block 106. The goal is to deliver messages consistently in an order that all hosts agree to, without requiring sources to know the constituency of the receiver set. Multicast group information is assumed to be available from a session directory service.

To allow selective addressing of hosts and dynamic election of an ordering node, ON, the TOM protocol introduces a labeling mechanism recently proposed for reliable multicast in the tree-based protocol Lorax (see, e.g., B. N. Levine et al., "The case for reliable concurrent multicasting using shared ack trees", Proc. ACM Multimedia, pages 365-376, Nov. 1996), and for multicast routing. Labels allow for open-ordered multicast, such as the addressing of specific nodes in the tree without the need to manifest a

separate multicast group or to reveal IP-addresses, wherein self-routing of messages to their destinations is facilitated based on prefix comparisons. Each node  $i$  in the acknowledgment-tree is labeled with a unique label  $l(i)$ , which is the prefix of all children of  $i$ . The label alphabet is preferably implemented with a set of symbols having a defined order, such as integers or letters with lexicographic order, with the alphabet cardinality corresponding to the tree branching factor  $B$ . The heuristics to select an ordering node, ON, is as follows: for each set of messages destined to a particular multicast group, or set of hosts, an ordering node, ON, is elected, such as by virtue of being the node whose label is the longest common prefix among all node labels in the receiver set. Each ordering node, ON, gathers sequence number bids set *en route* by primary nodes, PNs, deciding on a globally valid number, and multicasts the respective message to the receiver set with a final and binding sequence number directive.

FIG. 3 illustrates the mechanics of the TOM protocol exemplified on a multinode tree 200. Node  $r$ , as the root of the tree, carries label  $l$ . Node  $d$  is the only child in this multicast session which carries the prefix of its parent  $r$ , concatenated with its own index of "0". All three sources of messages, nodes  $x$ ,  $y$ , and  $z$ , have labels of length five, being positioned at depth five in the tree. An important principle in using labels for the ordering procedure is to create a confluence of messages at strategically optimal nodes in the tree for ordering a number of messages arriving in the same time window. Rather than depending on a statically assigned ordering node, the ordering node, ON, is



dynamically-selected per transmission, preferably as the node having the longest common prefix among the sources of pending messages in the targeted multicast group, without the need to pass an election token among nodes.

Consider the case that nodes  $x$ ,  $y$ , and  $z$  have messages to be multicast to a multicast group  $Rec = \{x, y, z, a, b, c, d, e, f\}$ . Each source multicasts its message to  $Rec$ , where it is entered in the order of collective arrival into  $uw$ . Control messages  $m_x^h$  and  $m_y^h$  are routed from source nodes, SNs,  $x$  and  $y$ , respectively, across their parents to the first common prefix node  $c$ , which are intermittently ordered at  $c$  with revised sequence numbers, and percolated up in the tree to node  $d$ , where message header  $m_x^h$  is also arriving. At any node on the path, a bitmask operation on the matching prefix indicates which messages must be up-routed, or handled locally. At node  $d$  it is determined that its label "10" matches the longest common prefix of SN labels  $l(x)$ ,  $l(y)$ ,  $l(z)$ . Hence, ordering nodes, ONs,  $(m_x, m_y, m_z) = d$  wherein node  $d$  sequences and multicasts the updated message headers to  $Rec$  to signal that the associated messages can be delivered. Once each receiver in  $Rec$  receives the ordering information per message  $m$  with  $of = true$  from the ordering node, ON, it shifts  $m$  into the  $ow$ , where the heading element is first delivered to end-processes.

Similarly, messages to a multicast group located in a left subbranch of the acknowledgment tree can be handled locally by the ordering node, ON, of that group, without affecting any nodes in other segments of the tree. The only overhead incurred in the ordering process is the control message unicast from source nodes, SNs, to

some ordering node, ON, plus one multicast to the receiver set. Total order is hence achieved in a diffusing computation, wherein the ordering process is carried out along with the message multicast, however, neither are receiver nodes, RNs, burdened with sorting out the messages, and they do not require knowledge of the identity of the  
5 ordering node, ON. Through the percolation process from source node, SN, to ordering node, ON, usage of the same sequence number for a specific message to all receivers in a multicast group is guaranteed.

10  
15  
20  
25  
30  
35  
40  
45  
50  
55  
60  
65  
70  
75  
80  
85  
90  
95  
100  
105  
110  
115  
120  
125  
130  
135  
140  
145  
150  
155  
160  
165  
170  
175  
180  
185  
190  
195  
200  
205  
210  
215  
220  
225  
230  
235  
240  
245  
250  
255  
260  
265  
270  
275  
280  
285  
290  
295  
300  
305  
310  
315  
320  
325  
330  
335  
340  
345  
350  
355  
360  
365  
370  
375  
380  
385  
390  
395  
400  
405  
410  
415  
420  
425  
430  
435  
440  
445  
450  
455  
460  
465  
470  
475  
480  
485  
490  
495  
500  
505  
510  
515  
520  
525  
530  
535  
540  
545  
550  
555  
560  
565  
570  
575  
580  
585  
590  
595  
600  
605  
610  
615  
620  
625  
630  
635  
640  
645  
650  
655  
660  
665  
670  
675  
680  
685  
690  
695  
700  
705  
710  
715  
720  
725  
730  
735  
740  
745  
750  
755  
760  
765  
770  
775  
780  
785  
790  
795  
800  
805  
810  
815  
820  
825  
830  
835  
840  
845  
850  
855  
860  
865  
870  
875  
880  
885  
890  
895  
900  
905  
910  
915  
920  
925  
930  
935  
940  
945  
950  
955  
960  
965  
970  
975  
980  
985  
990  
995

Labels allow open ordered multicast, such as the addressing of specific nodes in the tree with an ordered message sequence without the need to manifest a separate multicast group, and for self-routing of messages to their destinations based on prefix comparison. FIG. 4 sets forth an embodiment of the ordering algorithm 300 of TOM( ) that an ontree host  $i$  may utilize to send a message  $m$  totally ordered to a receiver set  $Rec$ , wherein hosts are assumed to carry prefix labels. Procedure  $TOM\_send()$  multicasts a message to the receiver set and unicasts the control header towards the dynamically elected ON;  $TOM\_cast()$  self-routes messages to a receiver based on prefix labels; and  $TOM\_receive()$  checks, whether a node is EN, PN, ON, or RN and takes action accordingly.

Consider the special case of ordering with this mechanism, in response to messages which are to be sent to two different, but overlapping, multicast groups. An  
20 example of the overlapping groups are  $G_1 = \{a, b, c\}$  and  $G_2 = \{c, d, e, f\}$  wherein  $G_1 \cap G_2 = c$ . Nodes in each group must receive a given message sequence in total

order, and node  $c$  shall not receive contradictorily ordered messages. This situation can be resolved, if individual membership within the target groups is known. Instead of choosing the node with the longest common prefix as the ordering node, ON, the nodes with multiple membership become the ordering cores for a transmission, and prescribe their sequencing decisions to their respective ordering node, ON. In the present case, node  $c$  will be instrumental in informing node  $d$  about the sequence in group  $G_1$ , such that node  $d$  can thereby construct a sequence compatible with  $G_2$ .

While total ordering of messages within one or more destination multicast groups is ensured, causal order among messages is not preserved in the above algorithm. To provide causality, the sequence numbers of messages to be ordered must incorporate encoded causal dependency information before reaching the ordering node, ON. By way of example, the encoding of causality information may be achieved by utilizing Lamport clocks which are maintained by all nodes belonging to a multicast group, and updating sequence numbers in the staggered ordering process to preserve the causal relations. To implement atomicity in delivery, that is, either all receiver nodes, RN, within  $Rec(m)$  will receive message  $m$ , or no message at all. Another message exchange must be introduced between receiver nodes, RNs, and ordering nodes, ONs, such that all receiver nodes, RNs, signal their reception of  $m$  and  $m^h$  to the ordering node, ON, and the ordering node, ON, is required to send another  $ok\_to\_deliver(m)$  signal for the receiver node, RN, to collectively proceed with delivery.

Resilience is another important aspect in TOM operation that is now briefly discussed. Ordering can be linked with several types of reliability, including (1) no guarantees on reliability of ordered deliveries, (2) the assumption of only inconsistent deliveries with failed hosts, (3) inciting roll-backs at operational hosts to repair inconsistent deliveries, and (4) the assumption that inconsistencies do not occur. Furthermore, another set of choices address the requirement to deliver a message, and the recipients to which the delivery guarantee is to be extended. In the event of host or link failures, the ordering tree may be partitioned into subtrees, each of which may continue to run TOM. The disappearance of an ordering node, ON, will be preferably remedied by replacement with the next common node in the destination set according to the label semantics. In operational subgroups, the semantics of reliable delivery is preserved for all multicast operations. Failure and recovery events must be made known to all operational hosts in an ordered fashion. Partitioned subbranches of the ordering tree may rejoin as soon as communication paths between them are reestablished. A link failure is detected, when a host fails to probe a neighbor node on the tree before expiration of a local timer. A host failure is detected, when a host with a pending queue of messages does not receive an expected message within a given timeout period.

#### 4. Taxonomy and Performance Comparison

Predominant ordering paradigms are classified using reliable broadcast or multicast into two main classes, as depicted in FIG. 5, wherein (1) *geometry-independent* protocols include *symmetric*, *two-phase*, and *centralized* solutions; while

(2) *geometry-dependent* protocols include *ring-based* and *tree-based* solutions. The following describes these paradigms and analyzes performance metrics to provide a performance analysis with the TOM protocol which operates on geometry-dependent tree-based protocols.

5           A number of multicasting schemes may involve all hosts in the ordering process in a decentralized way, using message stability properties, in contrast to solutions that burden one or a few of the hosts with the responsibility to order messages on behalf of the hosts in a multicast group. The main problem in the first case is to reach consensus among hosts on ordering patterns, the problem in the second case is to elect sequencer nodes. The present taxonomy contrasts the distinction between symmetric and token-site algorithms proposed by Rodriguez et. al. ("Totally ordered multicast in large-scale systems", Proc. of the 16<sup>th</sup> Int. Conf. on Distributed Computing Systems, pages 503-410, May 1996), which only accommodates symmetric protocols utilizing token-passing methods, and does not provide for tree-based ordering.

15           The processing of load  $X$  is evaluated at involved hosts and the message overhead  $M$  required to successfully multicast a message, in order, from a source node to all receivers. IP-multicast is assumed as the dissemination model for all schemes, although all schemes except TOM have been proposed in broadcast systems. The goal of this comparison is not an elaborate modeling of the many  
20 possible nuances and optimizations of ordering schemes in conjunction with reliable multicast, but rather a plain comparison of the fundamental working structure of ordering solutions. To this end, the evaluation does not include loss probabilities and assumes

that all schemes consistently use *sender-initiated* or *receiver-initiated* error recovery. Sender-initiated models place the burden for processing acknowledgments and requests for corrupt or lost packets on the transmission source, opposite to receiver-initiated solutions, wherein the retransmissions are performed in local groups among receivers and sources that are contacted only in the case of unrecoverable packet-loss. It should be appreciated that receiver-initiated protocols achieve improved scalability, largely due to the fact that sources are generally contacted only in the case of packet loss.

The notation used is as follows:  $s$  is the number of sources transmitting a message  $m$  destined for the same receiver, or receivers, at any given time, wherein each sender is assumed to also be a receiver;  $r$  is the number of receivers of message  $m$  in the receiver set  $Rec(m)$ ;  $X_f$  is the time required to feed a packet from a higher protocol layer;  $X_p$  is the time to process the transmission of a packet, including the time required for retransmissions;  $X_{\#}$  is the time to process a sequence number check;  $Y_p$  is the time to process a newly received packet;  $Y_f$  is the time to deliver a packet to an end process;  $X^w$  is the processing overhead per message in protocol  $w = \{S, 2P, C, R, T^{MP}, T^{MG}, T^{TOM}\}$ .  $M$  represents the number of transmissions required for all receivers to receive a message in a given order.

#### 4.1 Geometry-Independent Protocols

Reliable broadcast solutions are largely designed for fault-tolerant, asynchronous, distributed systems which utilize protocols that are geometry-

independent, for example wherein all hosts are assumed to be fully connected with one other, and wherein the routing between hosts does not presume any prearranged host geometry. Symmetric, two-phase, and centralized solutions are subsumed under this geometry-independent paradigm. Centralized ordering may also be classified as a star-  
 5 geometry, but the central node is typically chosen from all the nodes in an *ad-hoc* manner based on a predetermined election or token-passing scheme.

#### 4.1.1. Symmetric Ordering

In *symmetric* ordering schemes ( $S$ ), all hosts participate in the ordering process in a decentralized manner, analogous to a voting process, using message stability properties. A source node (SN) disseminates messages reliably to all hosts, which assigns a timestamp to each message and places it in a pending buffer; for each message  $m$ . Participant hosts (SN and RN) agree on a unique order number using timestamp information by running a consensus protocol. Messages with an assigned order number are shifted to the delivery queue and delivered to end processes in the globally binding order. It will be appreciated, therefore, that the number of messages to  
 5 be exchanged is a function of the number of hosts within the system that are involved in the ordering process. With  $X_C$  denoting the extra cost for the consensus protocol, the expected overhead of a generic symmetric protocol at the source node (SN) and receiver node (RN) is given by:

$$\begin{aligned}
 X_{SN}^S &= X_f + rX_p \\
 X_{RN}^S &= s(Y_p + X_{\#} + rX_c + Y_p)
 \end{aligned}
 \tag{1}$$

Utilizing broadcast communication, a source node sends a message to  $r - 1$  receivers, which in turn send  $r - 1$  messages to agree on the final sequence number, wherein  $M_{BC} = s((r - 1) + r(r - 1))$ , that is  $O(sr^2)$  for  $s$  sources. With multicast and  $r < n$  receivers,  $M = s(1 + 2r)$ , that is one multicast message to all receivers, one  
 5 multicast per each of the  $r$  receivers to each other, and one timestamp sweep from all receivers to the source. Protocols with fault-tolerance measures may incur significantly higher cost factors.

#### 4.1.2. Two Phase Ordering

Four communication steps are required when utilizing *two-phase* ordering ( $2P$ ).  
 A source sends a message  $m$  to a multicast group, whereupon each receiver assigns a  
 10 priority number to the message, places  $m$  as pending in its local queue, and returns the priority number to the source. The source selects the highest number and sends it to all receivers, thereby replacing the original number with the new one, tags the message as deliverable, reorders the queue, and delivers the messages at the head of the queue.  
 15 Expected overhead at the source node (SN) and the receiver node (RN) is given by:

$$X_{SN}^{2P} = X_f + r(Y_p + X_{\#} + 2X_p) \quad (2)$$

$$X_{RN}^{2P} = s(2Y_p + X_{\#} + X_p + Y_f)$$

If it is assumed  $r \geq s$ , then  $X^{2P} = \max(X_{SN}^{2P}, X_{RN}^{2P}) = O(r)$ . Given one  
 message multicast from  $s$  sources to  $r$  receivers, a number of control messages  $r$   
 20 with priority numbers are sent back to each source, while a final control message must



be multicast from the source to the receiver set for each message, such as  $M = s(1 + r)$ .

#### 4.1.3. Centralized Ordering

In *centralized* ordering ( $C$ ) a source node (SN) transmits a message  $m$  to a sequencer host, which assigns a unique number to  $m$ , and forwards it to the receiver set  $Rec(m)$ , where it is ultimately delivered to end-processes in the order prescribed by the sequence numbers. The sequencer role may rotate among hosts. The expected overhead at SN, ON, and RN is thereby given by:

$$X_{SN}^C = X_f + X_p \quad (3)$$

$$X_{ON}^C = s(Y_p + X_{\#} + rX_p)$$

$$X_{RN}^C = s(Y_p + Y_f)$$

Hence  $X^C = O(sr)$ , and  $M = s + r$ , consisting of  $s$  messages from sources to the ordering node (ON), and one multicast per message from ordering node (ON) to all receivers. If the source node (SN) is the same as the ordering node (ON), then one step is eliminated.

### 4.2 Geometry-Dependent Protocols

*Geometry-dependent* protocols presume a specific host topology to route ordering information.

#### 4.2.1 Ring-based Ordering

In *ring-based* ordering ( $R$ ) a logical ring imposes a transmission path between

hosts, wherein each host is only required to communicate with its predecessor and its successor in the ring. To multicast a message, a host must possess the token. The token contains requests for messages to be resent and the highest sequence number for any message broadcast on the ring. Each host maintains an input buffer containing pending messages with assigned sequence numbers. On receipt of the token, the host completes processing of the messages in its buffer by adjusting sequence numbers, resends messages requested in the token, updates the token information and forwards the token. Messages are sent to end processes when marked as deliverable. Each source node (SN), as a token-site, assumes the role of an ordering node (ON). With  $X_{tk}$  indicating the token transfer time, the expected overhead at the source node (SN) and the receiver node (RN) in a single ring is given by:

$$X_{SN}^R = X_f + X_p + r(Y_p + X_{\#} + X_p) + X_{tk} \quad (4)$$

$$X_{SN}^R = s(Y_p + X_{\#} + Y_f)$$

Hence  $X^R = O(r)$ , if  $r > s$ , and the minimum message overhead is given by  $M = 2n/k$ , where  $2n$  is the number of token transfers required to accept  $k$  multicast messages in a ring of  $n$  nodes. Assuming that  $k = 1$  with  $s$  sources, and despite  $r < n$  receivers,  $M = 2sn$ .

#### 4.2.2 Tree-based Ordering

For *tree-based* ordering ( $T$ ), the MP protocol and the metagroup approach (MG) are compared with TOM. It will be appreciated that these current tree-based reliable multicast protocols do not provide ordering. Common to MP, MG, and TOM is the

element of distributing the ordering responsibility and load across several nodes on the tree. The IMP and MG protocols utilize group membership information to cluster nodes for optimized message delivery, in contrast to which the TOM protocol utilizes the end-to-end multicast topology.

5           The MP protocol include two operating phases (1) the transmission from the source to a primary host, and (2) the transmission from this host to the receivers. MP builds a plethora of propagation trees, wherein hosts in the intersections of multicast groups are chosen as hop nodes, such as the roots of subtrees. A message is first sent to these primary hosts, and then propagated downward in the tree toward the receiver hosts, being ordered on their propagation path, and finally unicast to the receiver hosts.

10           The MG protocol clusters hosts from overlapping multicast groups into metagroups, which do not overlap. Each group has a primary metagroup (PM), and in each metagroup one member is assigned to be a manager. Metagroups are organized in a plurality of propagation trees, such that the PM of a group is the ancestor of all other metagroups of the same group in the tree. Messages destined to multicast group  $G$  are first sent to the primary node  $PM(G)$ , which propagates the messages along the tree to all other metagroups, which are subsets of  $G$ . The manager of a metagroup broadcasts a message to other members in its metagroup.

15           The drawback with the MP and MG protocols is the need to compute a logical propagation or metagroup tree per-source as overlays to the end-to-end geometry, which requires that in order to construct such a tree, the computation host, or hosts, must recognize the membership of all groups. This approach is operable only for

closed multicast and static groups, and the cost may be rationally amortized only for long-duration transmissions between hosts. The processing overhead common to all tree-based schemes is:

$$X_{SN}^T = X_f + X_p \quad (5)$$

$$X_{ON}^T = B(Y_p + X_{\#} + X_p)$$

$$X_{RN}^T = Y_p + Y_f$$

Hence generally  $X^T = O(B)$ , where  $B$  indicates the branching factor of the tree.

With multicast,  $M^{MP} = s(1 + d)$  messages are required, with one message from each of the  $s$  sources to the primary destination in the subtree, and one broadcast at each level of the subtree, where  $d$  is the subtree depth. The MG protocol has three operational phases and requires one message to  $PM(G)$ ,  $d$  messages to the managers of the deepest metagroups at depth  $d$  in the subtree, and another  $k$  messages to the members of the  $k$  metagroups containing the target multicast group, wherein  $M^{MG} = s(1 + d + k)$ .

It will be appreciated that TOM requires a multicast from  $s$  source nodes (SNs) to the receiver set, and  $p$  unicasts from the source node (SN) to the ordering node (ON), where  $p$  is the average path length, and one final multicast from the ordering node (ON) to the receiver node (RN), wherein  $M^{TOM} = s(2 + p)$ .

#### 4.3. Results and Comparisons

Table 1 summarizes expected message costs and delays for the described

protocols. Centralized and two-phase approaches incur only two, and three message exchange phases, respectively, however, the messaging is concentrated on specific hosts in the session which are subject to failure and bottlenecks. Rings engage all hosts in a session in the transmission process, even when a source and multicast receiver group constitute only a small portion of the entire session. Trees allow for selective engagement of hosts on those subbranches or local groups, which are actually affected by the message processing.

It is assumed that there are as many sources as receivers,  $r = n$  and  $s = 1$ . In the graph the cost to compute and maintain the propagation infrastructure is ignored, although the anticipated overhead for the MG and MP protocols is substantial in contrast with the TOM protocol which simply relies on a given acknowledgment tree. The session size is varied between  $n = [1, 1000]$ , with  $r = n/10$  as the average size of a receiver multicast group. The tree-depth of the MP protocol has been projected between  $d = [1, 8]$  for simulations with  $n = 200$  and average group size  $g = [5, 40]$ . The tree depth for a metagroup tree has been projected between  $d = [1, 5]$  for up to 40 metagroups with  $g = 50$ , and an overlapping degree of 10. It is also assumed that each source sends only one multicast message per transmission cycle. Simulations for the Lorax protocol have indicated that optimal ACK trees are built when each node supports at least  $B = 5$  neighbors. To provide a baseline comparison, the average depth of a subbranch in a tree according to the MP and MG protocols is chosen as  $d = \log B^r$ , where  $B = 5$ , depicts the average node degree. The average path length

according to the TOM protocol is chosen as  $p = h/2$ , because roughly half of the height  $h$  of the tree needs to be traversed in converging on a particular ordering node (ON). It should be noted that a message comparison provides a limited view on the relative performance of the protocols, because parallelism in message processing, the processing overhead at various nodes, and the shape of the tree would need to be considered in a more precise way. However, concentrating on  $M$  alone is sufficient to express fundamental differences between the approaches. FIG. 6 plots the multicast message cost of the various schemes under given assumptions.

The results only represent performance of the discussed protocols under one particular scenario, namely genuine multicast utilizing a single transmission source. The multiple source case would reinforce that the throughput of a generic tree-based protocol for ordered reliable multicast scales better with receiver set due to locus and execution of sequencing. Symmetric methods exhibit the least amount of scalability, as a result of requiring that all nodes be involved in processing messages from all other nodes. If all nodes broadcast at the same time, latency may be low, but a consensus protocol must be run. Two-phase, centralized, and ring solutions have similar message overhead. The use of the ring solutions, however, may permit higher-concurrency, although a drawback arises for large sessions due to latency increases. The centralized ordering method is reasonably efficient when limited to a few hosts, however, it is subject to potential bottlenecks and results in a single point of failure, which is particularly risky when utilized for large sessions. A logical hop between hosts within the MP and MG protocols may require multiple hops across long distances in the

multicast routing tree, in contrast to the TOM protocol, which operates under the assumption that the structure of the ACK-tree mirrors the path information in the multicast routing tree, rather than using separate propagation graphs. Comparing the three tree-based methods, it will be appreciated that the TOM protocol of the present invention provides equal, or improved, performance in relation to either the MG or MP protocols. TOM also spreads the computational load of ordering packets over multiple nodes in the tree, and is well suited for dynamically altering multicast groups, rather than catering to static membership and long-lived transmissions.

## 5. Conclusions

The present invention provides for the addition of ordering services to tree-based concurrent reliable packet multicasting which is essential to a growing number of Internet applications supporting telepresence and near-synchronous information sharing. Considering the use of reliable multicasting for these applications, it has been observed that ordering services have not been integrated as a component in the currently available data dissemination methods. The TOM protocol, however, stands in contrast to previous reliable broadcast solutions tailored to local area networks, wherein ordering was performed assuming symmetric communication, centralized, ring-based, or propagation graph schemes. It will be appreciated that the TOM protocol is readily applicable to a number of multicasting applications. Furthermore, although TOM is directed towards the addition of an ordering capability for use within reliable concurrent multicasting, such as defined by Lorax, it may be equally deployed in other frameworks, for example, TMTP with domain managers, and in RMTP with designated receivers as

intermediate ordering nodes.

Accordingly, it will be seen that the TOM protocol is solution directed at providing reliable multicast trees, using staggered ordering of messages on their paths from sources to receivers. The workload of executing the ordering protocol when utilizing the TOM protocol is distributed across the nodes wherein the infrastructure being utilized for packet ordering is cohesive and results in reliable operation. The addition of address labels yields efficient ordering for multiple groups and subgroups. In contrast with other prominent solutions, the TOM protocol does not require computation of separate graphs for propagating ordering information. The TOM multicast ordering protocol implements ordering in a diffusing computation, wherein messages are ordered on their delivery paths from sources to receivers, and each node communicates only with its children and parent node instead of the entire multicast group. A taxonomy has been proposed for ordering schemes integrating reliable broadcast and multicast solutions. A simple performance comparison has illustrated that ordering within trees surpasses the use of contending solutions in terms of scalability, efficiency, and practicality. It should be appreciated that although the description of distributed multicasting solution for tree-based multicasting was exemplified with method steps and pseudocode procedures, it may be implemented with numerous variations by one of ordinary skill in the art without departing from the teachings of the present invention.

Although the description above contains many specificities, these should not be construed as limiting the scope of the invention but as merely providing illustrations of some of the presently preferred embodiments of this invention. Therefore, it will be



